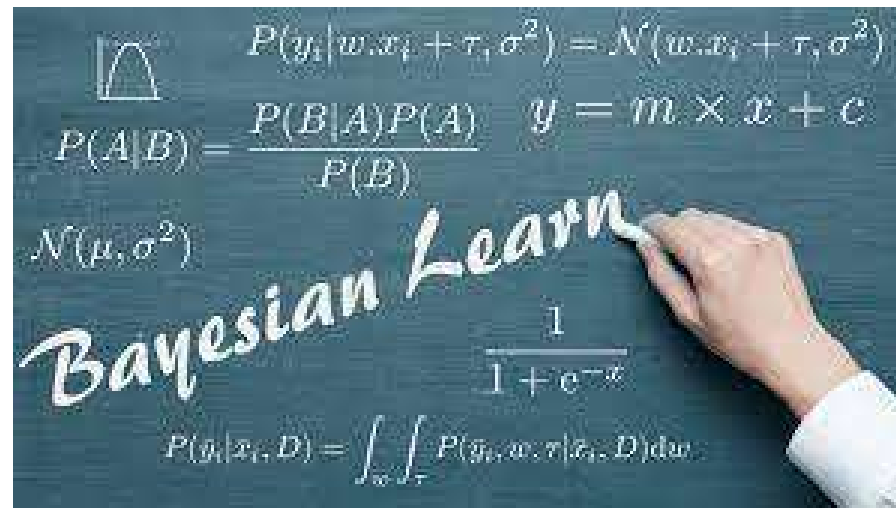


CT-562 MACHINE LEARNING

NED University of Engineering & Technology



BAYESIAN LEARNING

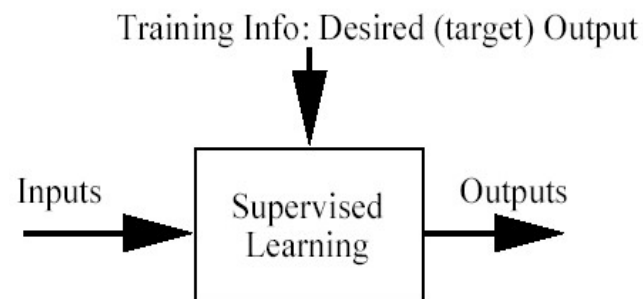
OVERVIEW

Bayes rule & turn this into a classifier

- E.g. How to decide if a patient is ill or healthy, based on
 - A probabilistic model of the observed data
 - Prior knowledge

CLASSIFICATION PROBLEM

- Training data: examples of the form $(d, h(d))$
 - where d are the data objects to classify (inputs)
 - and $h(d)$ are the correct class info for d , $h(d) \in \{1, \dots, K\}$
- Goal: given d_{new} , provide $h(d_{\text{new}})$



Error = (target output - actual output)

A WORD ABOUT THE BAYESIAN FRAMEWORK

- Allows us to combine observed data and prior knowledge
- Provides practical learning algorithms
- It is a generative (model based) approach, which offers a useful conceptual framework
 - This means that any kind of objects (e.g. time series, trees, etc.) can be classified, based on a probabilistic model specification

BAYES' RULE

$$p(h | d) = \frac{P(d | h)P(h)}{P(d)}$$

Who is who in Bayes' rule

$P(h)$: prior belief (probability of hypothesis h before seeing any data)

$P(d | h)$: likelihood (probability of the data if the hypothesis h is true)

$P(d) = \sum_h P(d | h)P(h)$: data evidence (marginal probability of the data)

$P(h | d)$: posterior (probability of hypothesis h after having seen the data d)

Understanding Bayes' rule

d = data

h = hypothesis

Proof. Just rearrange:

$$p(h | d)P(d) = P(d | h)P(h)$$

$$P(d, h) = P(d, h)$$

the same joint probability

on both sides

PROBABILITIES

- Have two dice h_1 and h_2
- The probability of rolling an i given die h_1 is denoted $P(i|h_1)$. This is a conditional probability
- Pick a die at random with probability $P(h_j)$, $j=1$ or 2 . The probability for picking die h_j and rolling an i with it is called joint probability and is $P(i, h_j) = P(h_j)P(i|h_j)$.
- For any events X and Y , $P(X, Y) = P(X|Y)P(Y)$
- If we know $P(X, Y)$, then the so-called marginal probability $P(X)$ can be computed as
$$P(X) = \sum_Y P(X, Y)$$
- Probabilities sum to 1. Conditional probabilities sum to 1 **provided that their conditions are the same.**

DOES PATIENT HAVE CANCER OR NOT?

A patient takes a lab test and the result comes back positive. It is known that the test returns a correct positive result in only 98% of the cases and a correct negative result in only 97% of the cases. Furthermore, only 0.008 of the entire population has this disease.

1. What is the probability that this patient has cancer?
2. What is the probability that he does not have cancer?
3. What is the diagnosis?

$hypothesis1: 'cancer'$
 $hypothesis2: '\neg cancer'$ } hypothesis space H

– data : '+'

$$1. P(cancer | +) = \frac{P(+ | cancer)P(cancer)}{P(+)} = \frac{\dots\dots\dots}{\dots\dots\dots} = \dots\dots\dots$$

$$P(+ | cancer) = 0.98$$

$$P(cancer) = 0.008$$

$$P(+) = P(+ | cancer)P(cancer) + P(+ | \neg cancer)P(\neg cancer)$$

$$= \dots\dots\dots$$

$$P(+ | \neg cancer) = 0.03$$

$$P(\neg cancer) = \dots\dots\dots$$

$$2. P(\neg cancer | +) = \dots\dots\dots$$

3. Diagnosis ??

CHOOSING HYPOTHESES

- *Maximum Likelihood* hypothesis:
- Generally we want the most probable hypothesis given training data. This is the *maximum a posteriori* hypothesis:
 - Useful observation: it does not depend on the denominator $P(d)$

$$h_{ML} = \arg \max_{h \in H} P(d | h)$$

$$h_{MAP} = \arg \max_{h \in H} P(h | d)$$

NOW WE COMPUTE THE DIAGNOSIS

- To find the Maximum Likelihood hypothesis, we evaluate $P(d|h)$ for the data d , which is the positive lab test and chose the hypothesis (diagnosis) that maximises it:

$$P(+ | cancer) = \dots\dots\dots$$

$$P(+ | \neg cancer) = \dots\dots\dots$$

$$\Rightarrow \text{Diagnosis} : h_{ML} = \dots\dots\dots$$

- To find the Maximum A Posteriori hypothesis, we evaluate $P(d|h)P(h)$ for the data d , which is the positive lab test and chose the hypothesis (diagnosis) that maximises it. This is the same as choosing the hypotheses gives the higher posterior probability.

$$P(+ | cancer)P(cancer) = \dots\dots\dots$$

$$P(+ | \neg cancer)P(\neg cancer) = \dots\dots\dots$$

$$\Rightarrow \text{Diagnosis} : h_{MAP} = \dots\dots\dots$$

THE NAÏVE BAYES CLASSIFIER

- What can we do if our data d has several attributes?
- Naïve Bayes assumption: Attributes that describe data instances are conditionally independent given the classification hypothesis

$$P(\mathbf{d} | h) = P(a_1, \dots, a_T | h) = \prod_t P(a_t | h)$$

- it is a simplifying assumption, obviously it may be violated in reality
 - in spite of that, it works well in practice
- The Bayesian classifier that uses the Naïve Bayes assumption and computes the MAP hypothesis is called Naïve Bayes classifier
- One of the most practical learning methods
- Successful applications:
 - Medical Diagnosis
 - Text classification

EXAMPLE. 'PLAY TENNIS' DATA

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

NAÏVE BAYES SOLUTION

Classify any new datum instance $\mathbf{x}=(a_1, \dots, a_T)$ as:

$$h_{Naive\ Bayes} = \arg \max_h P(h)P(\mathbf{x} | h) = \arg \max_h P(h) \prod_t P(a_t | h)$$

- To do this based on training examples, we need to estimate the parameters from the training examples:
 - For each target value (hypothesis) h
 $\hat{P}(h) := \text{estimate } P(h)$
 $\hat{P}(a_t | h) := \text{estimate } P(a_t | h)$
 - For each attribute value a_t of each datum instance

Based on the examples in the table, classify the following datum \mathbf{x} :

$\mathbf{x} = (\text{Outl} = \text{Sunny}, \text{Temp} = \text{Cool}, \text{Hum} = \text{High}, \text{Wind} = \text{strong})$

- That means: Play tennis or not?

$$h_{NB} = \arg \max_{h \in [\text{yes}, \text{no}]} P(h) P(\mathbf{x} | h) = \arg \max_{h \in [\text{yes}, \text{no}]} P(h) \prod_t P(a_t | h)$$

$$= \arg \max_{h \in [\text{yes}, \text{no}]} P(h) P(\text{Outlook} = \text{sunny} | h) P(\text{Temp} = \text{cool} | h) P(\text{Humidity} = \text{high} | h) P(\text{Wind} = \text{strong} | h)$$

- Working:

$$P(\text{PlayTennis} = \text{yes}) = 9 / 14 = 0.64$$

$$P(\text{PlayTennis} = \text{no}) = 5 / 14 = 0.36$$

$$P(\text{Wind} = \text{strong} | \text{PlayTennis} = \text{yes}) = 3 / 9 = 0.33$$

$$P(\text{Wind} = \text{strong} | \text{PlayTennis} = \text{no}) = 3 / 5 = 0.60$$

etc.

$$P(\text{yes}) P(\text{sunny} | \text{yes}) P(\text{cool} | \text{yes}) P(\text{high} | \text{yes}) P(\text{strong} | \text{yes}) = 0.0053$$

$$P(\text{no}) P(\text{sunny} | \text{no}) P(\text{cool} | \text{no}) P(\text{high} | \text{no}) P(\text{strong} | \text{no}) = \mathbf{0.0206}$$

$$\Rightarrow \text{answer} : \text{PlayTennis}(x) = \text{no}$$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Given all the previous patients that I have seen, below are their symptoms and diagnosis:

chills	runny nose	headache	fever	flu?
Y	N	Mild	Y	N
Y	Y	No	N	Y
Y	N	Strong	Y	Y
N	Y	Mild	Y	Y
N	N	No	N	N
N	Y	Strong	Y	Y
N	Y	Strong	N	N
Y	Y	Mild	Y	Y

Do I believe that the patient with the following symptoms has a flu using Naïve Bayes Algorithm?

Chills	Runny nose	Headache	Fever	Flu?
Y	N	Mild	N	?

REMEMBER

- Bayes' rule can be turned into a classifier
- Maximum A Posteriori (MAP) hypothesis estimation incorporates prior knowledge; Max Likelihood doesn't
- Naive Bayes Classifier is a simple but effective Bayesian classifier for vector data (i.e. data with several attributes) that assumes that attributes are independent given the class.
- Bayesian classification is a generative approach to classification



THANK YOU